



Saturdays.AI



Saturdays.AI
Donostia



#2 Supervised Learning

by Saturdays.AI

Saturdays.AI
Machine Learning

Schedule

State of the course

Session 2 Review

Challenge

Notebook + resources

State of the course

#1 Cleaning & Exploratory Data Analysis ✓

#2 Supervised Learning ● Today!

#3 Decision Trees & Random Forest ➔ SOON

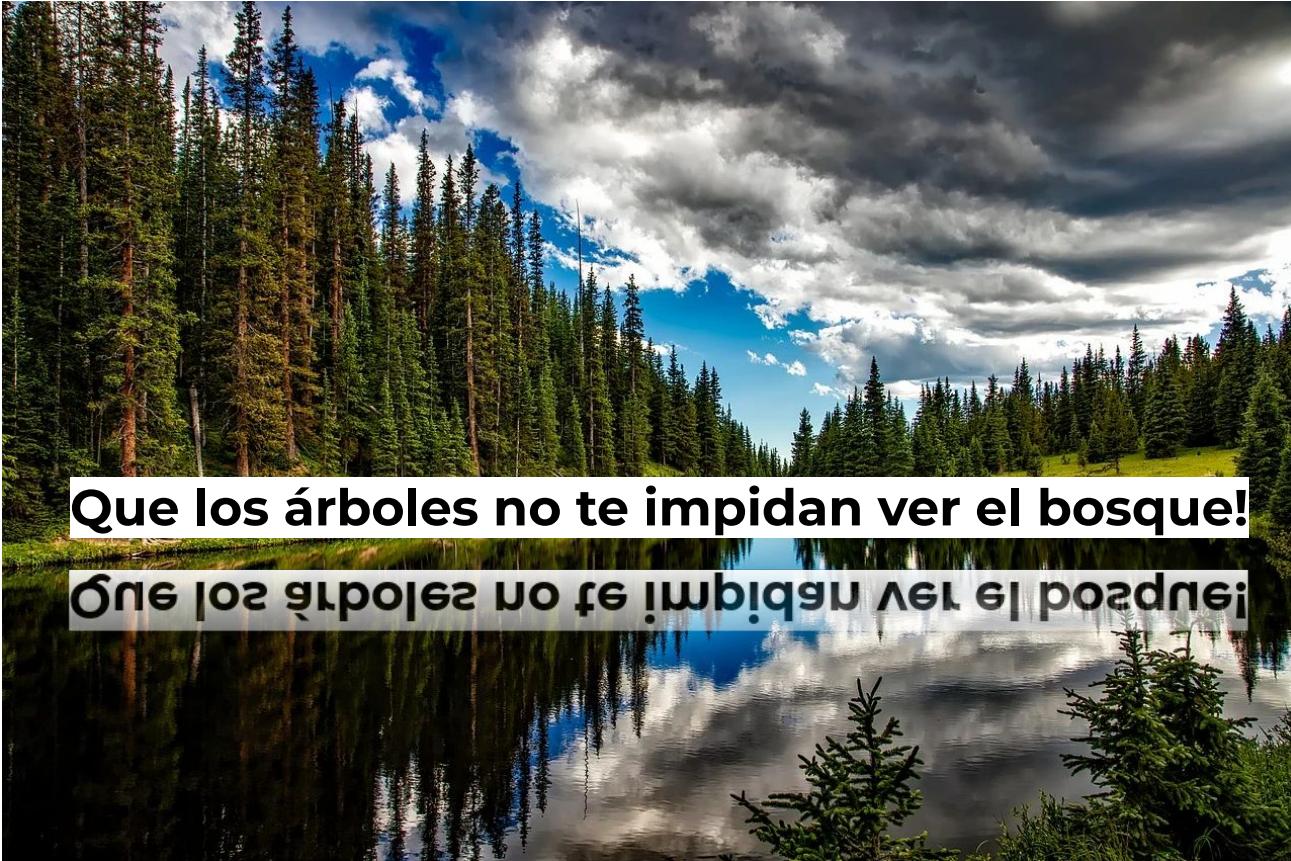
#4 Unsupervised Learning + Clustering ➔ SOON

#5 Time Series Analysis + Data Viz ➔ SOON

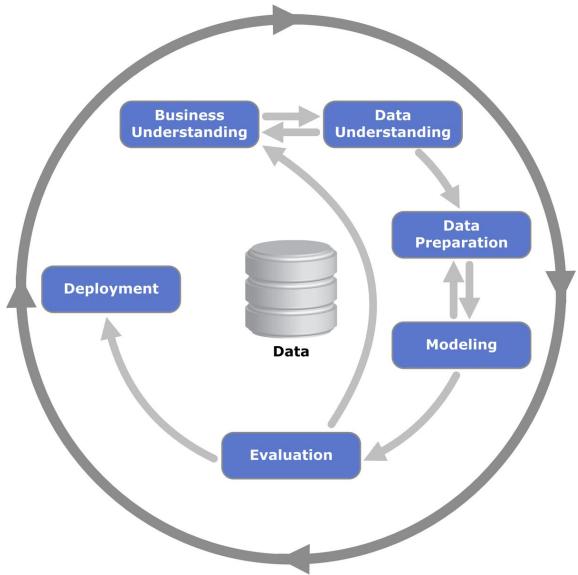
#6 Neural Networks, Gradient Descent ➔ SOON

#7 NLP ➔ SOON

Little general view...



Little general view... CRISP-DM



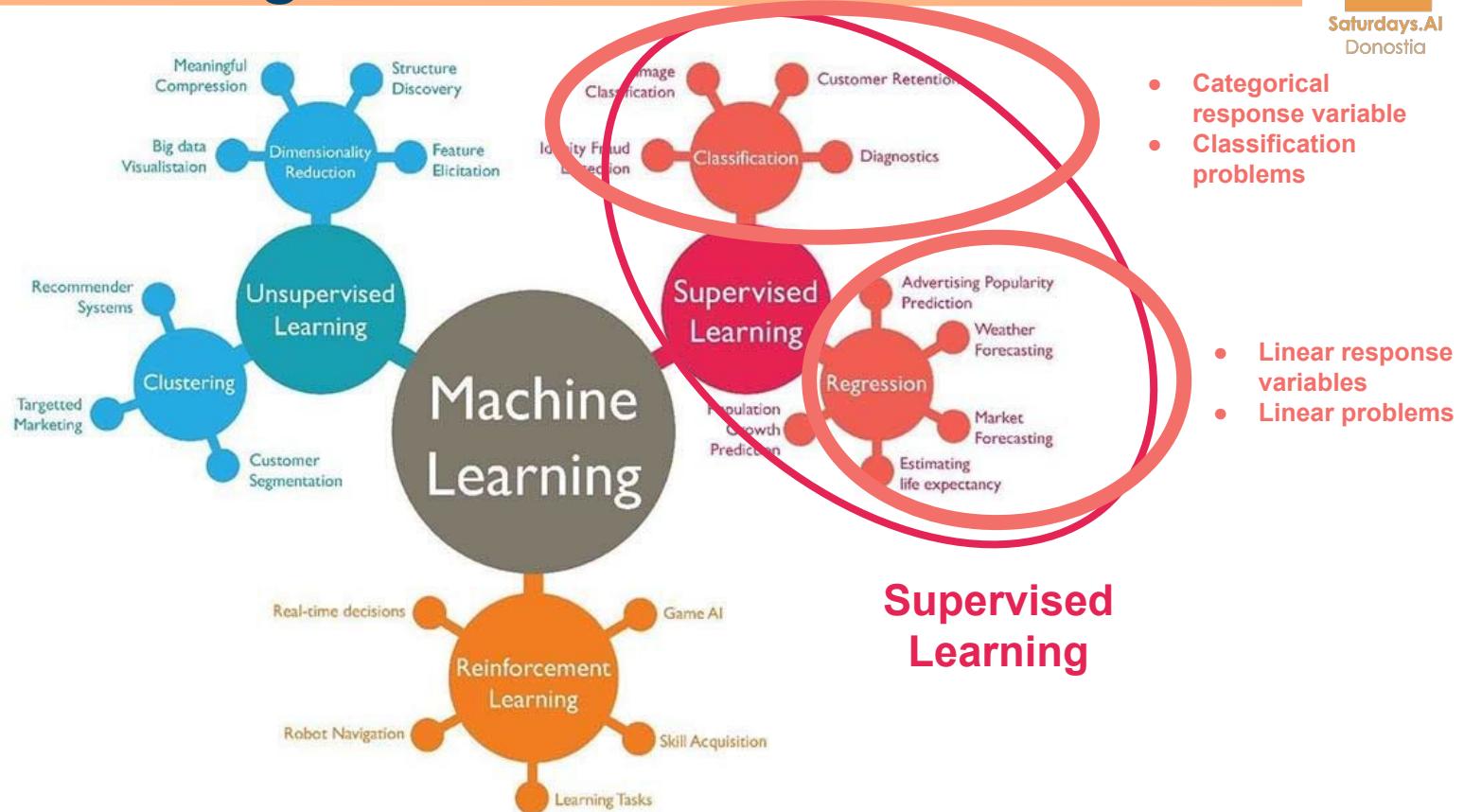
| Business Understanding | Data Understanding | Data Preparation | Modeling | Evaluation | Deployment |
|--|---|---|--|---|--|
| Determine Business Objectives <i>Background, Business Objectives, Business Success Criteria</i> Assess Situation <i>Inventory of Resources Requirements, Assumptions, and Constraints, Risks and Contingencies, Terminology, Costs and Benefits</i> Determine Data Mining Goals <i>Data Mining Goals, Data Mining Success Criteria</i> Produce Project Plan <i>Project Plan, Initial Assessment of Tools and techniques</i> | Collect Initial Data <i>Initial Data Collection Report</i> Describe Data <i>Data Description Report</i> Explore Data <i>Data Exploration Report</i> Verify Data Quality <i>Data Quality Report</i> | Select Data <i>Rationale for Inclusion/Exclusion</i> Clean Data <i>Data Cleaning Report</i> Construct Data <i>Derived Attributes, Generated Records</i> Integrate Data <i>Merged Data</i> Format Data <i>Reformatted Data</i> Dataset <i>Dataset Description</i> | Select Modeling Techniques <i>Modeling Technique, Modeling Assumptions</i> Generate Test Design <i>Test Design</i> Build Model <i>Parameter Settings, Models, Model Descriptions</i> Assess Model <i>Model Assessment, Revised Parameter Settings</i> | Evaluate Results <i>Assessment of Data Mining Results w.r.t. Business Success Criteria, Approved Models</i> Review Process <i>Review of Process</i> Determine Next Steps <i>List of Possible Actions Decision</i> | Plan Deployment <i>Deployment Plan</i> Plan Monitoring and Maintenance <i>Monitoring and Maintenance Plan</i> Produce Final Report <i>Final Report, Final Presentation</i> Review Project <i>Experience Documentation</i> |

Regression & Model Validation

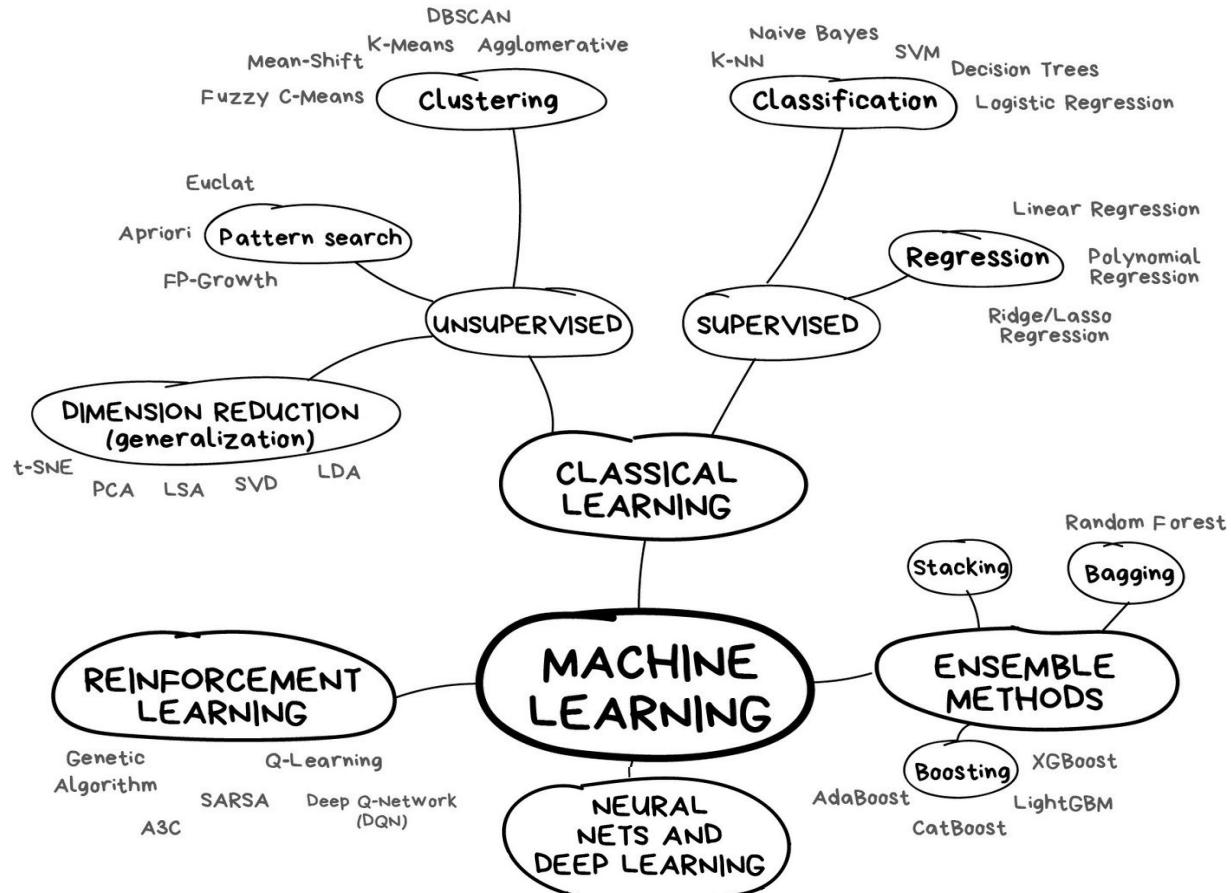
- General model topics
 - Overfitting/Underfitting
 - Bias/Variance
 - Cross-validation
 - Regularization
 - Model evaluation
- Models: Linear Models, Naive Bayes, AR, SVM

Are
You
Ready?

Machine learning



ML Algorithms



Supervised learning: linear problem (Ejemplo)

HOUSE PRICE Dataset



| Id | SalePrice | MSZoning | YearBuilt | LotArea | Street |
|----|-----------|----------|-----------|---------|--------|
| 1 | 208500 | RL | 2003 | 8450 | Pave |
| 2 | 181500 | RL | 1976 | 9600 | Pave |
| 3 | 223500 | RL | 2001 | 11250 | Pave |
| 4 | 140000 | RL | 1915 | 9550 | Pave |
| 5 | 250000 | RL | 2000 | 14260 | Pave |
| 6 | 143000 | RL | 1993 | 14115 | Pave |
| 7 | 307000 | RL | 2004 | 10084 | Pave |
| 8 | 200000 | RL | 1973 | 10382 | Pave |
| 9 | 129900 | RM | 1931 | 6120 | Pave |

Regression & Model Validation

Modeling
Phase

All models have 3 key components: a task, a learning methodology and a performance measure

Task

What is the problem we want our model to solve?

Learning
Methodology

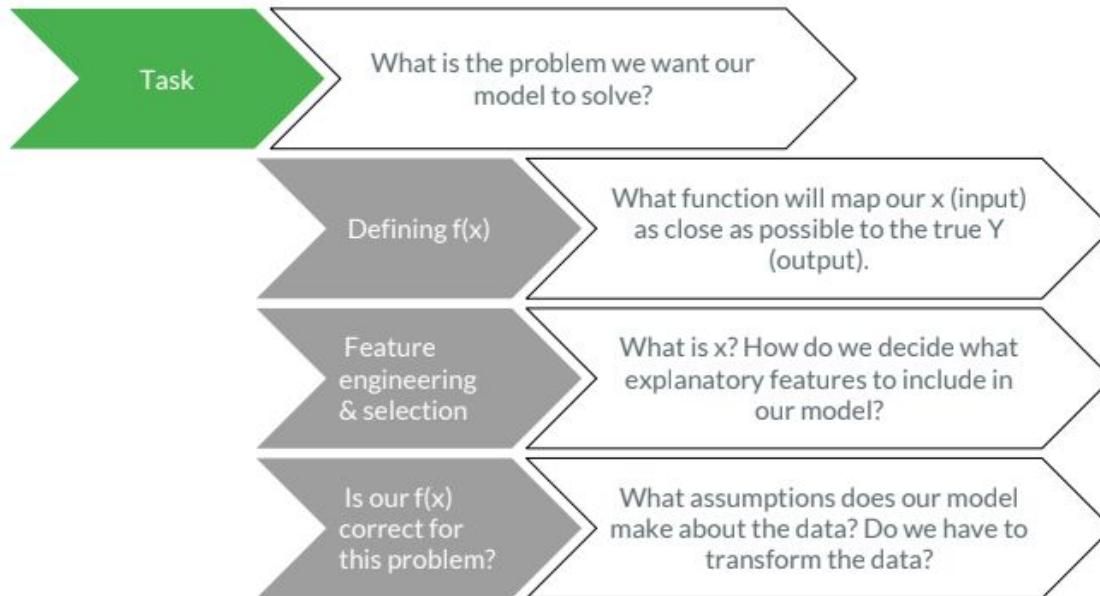
ML algorithms can be supervised or unsupervised. This determines the learning methodology.

Performance
Measure

Quantitative measure we use to evaluate the model's performance.

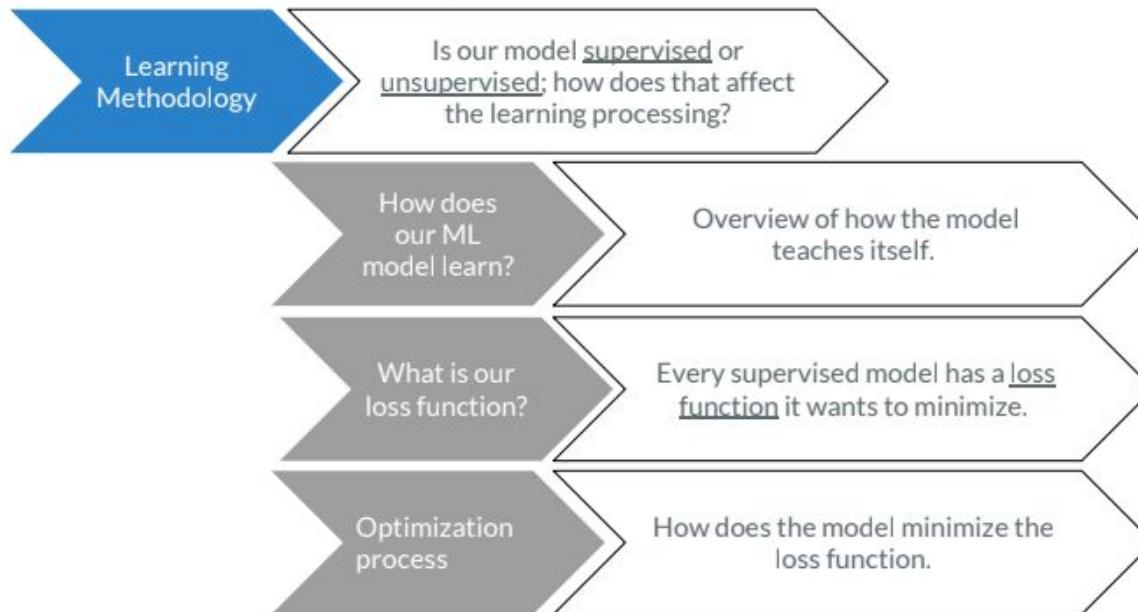
Regression & Model Validation

Today we are looking closer at each component of the framework:



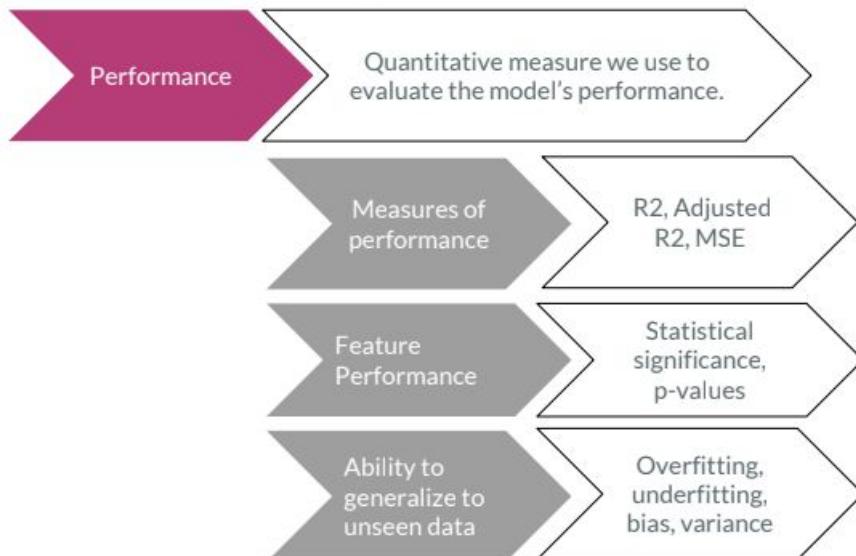
Regression & Model Validation

Learning methodology: how does the model learn the function that best maps x to the true y ?



Regression & Model Validation

Performance: How do we evaluate how useful the model is, and how we can improve it?



Regression Loss Function

Regression

Mean Square
Error/
Quadratic Loss

Mean Absolute
Error

Huber Loss/
Smooth Mean
Absolute Error

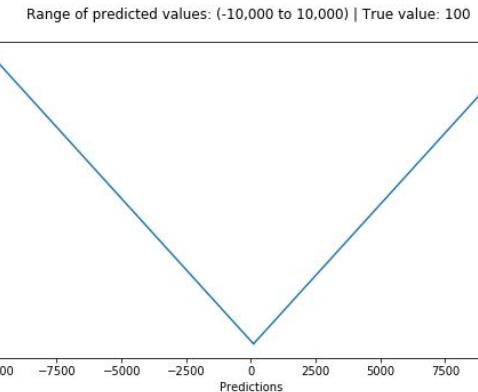
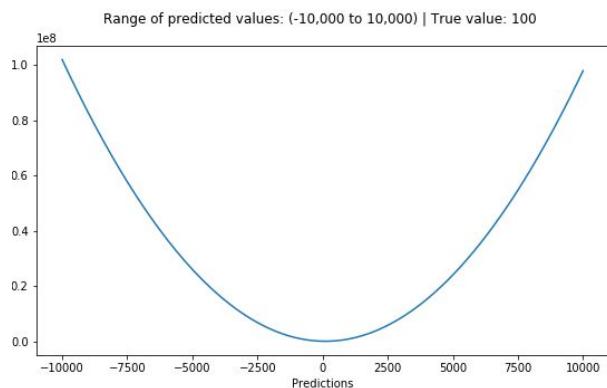
Log cosh Loss

Quantile Loss

Q: ¿Cómo de bien o mal vamos en el entrenamiento?

$$MSE = \frac{\sum_{i=1}^n (y_i - y_i^p)^2}{n}$$

$$MAE = \frac{\sum_{i=1}^n |y_i - y_i^p|}{n}$$



Regression Loss Function

Regression

Mean Square Error/
Quadratic Loss

Mean Absolute Error

Huber Loss/
Smooth Mean Absolute Error

Log cosh Loss

Quantile Loss

MAE vs. RMSE for cases with slight variance in data

| ID | Error | Error | Error ² |
|----|-------|-------|--------------------|
| 1 | 0 | 0 | 0 |
| 2 | 1 | 1 | 1 |
| 3 | -2 | 2 | 4 |
| 4 | -0.5 | 0.5 | 0.25 |
| 5 | 1.5 | 1.5 | 2.25 |

MAE: 1 RMSE: 1.22

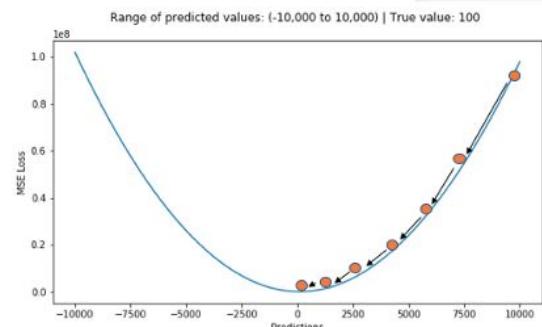
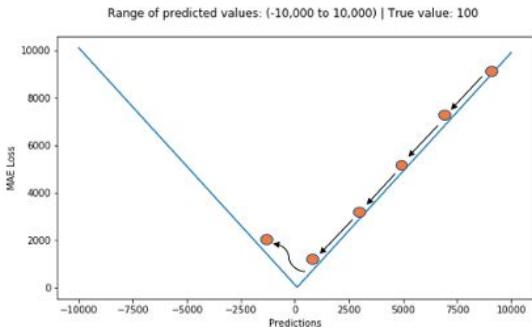
MAE vs. RMSE for cases with outliers in data

| ID | Error | Error | Error ² |
|----|-------|-------|--------------------|
| 1 | 0 | 0 | 0 |
| 2 | 1 | 1 | 1 |
| 3 | 1 | 1 | 1 |
| 4 | -2 | 2 | 4 |
| 5 | 15 | 15 | 225 |

MAE: 3.8 RMSE: 6.79

outlier

En la práctica, esto sucede en la computadora (cuando aprende estamos entrenando el modelo)...



Measures of performance

Todo parte de la MATRIZ DE CONFUSIÓN...

| | | Modelo | |
|----------|---------|------------|------------|
| | | PREDIJO SI | PREDIJO NO |
| Realidad | REAL SI | 80 | 10 |
| | REAL NO | 10 | 100 |

$$\text{Accuracy} = \frac{\text{TrueNegatives} + \text{TruePositive}}{\text{TruePositive} + \text{FalsePositive} + \text{TrueNegative} + \text{FalseNegative}}$$

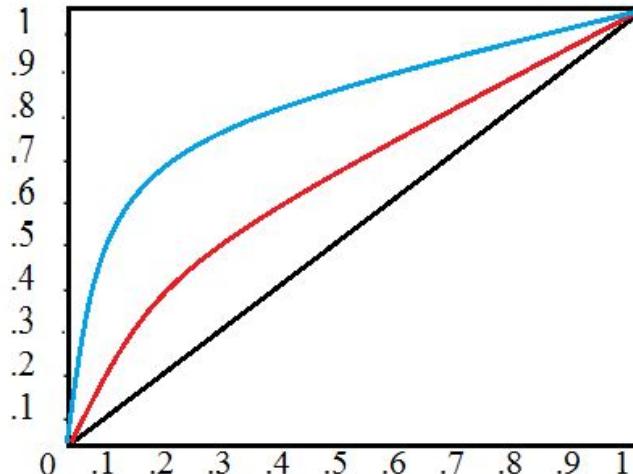
$$\text{recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

$$\text{precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

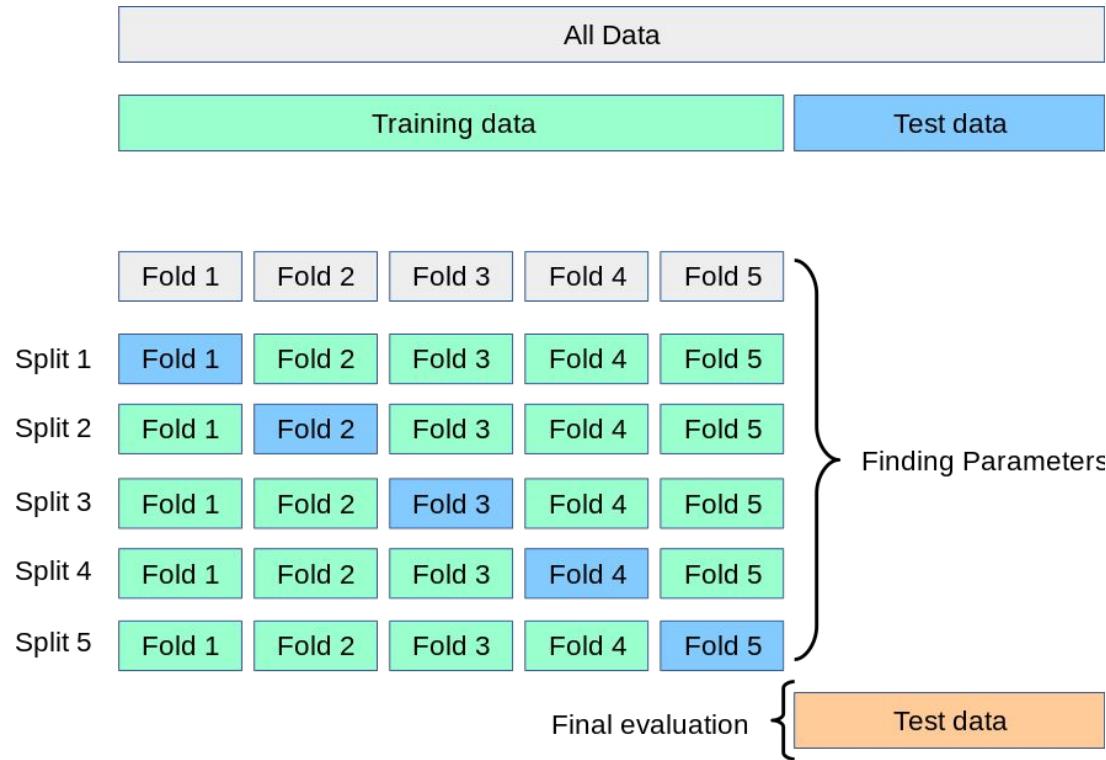
Measures of performance

ROC curve → Plots the true positive rate (TPR) versus the false positive rate (FPR) as a function of the model's **threshold** for classifying a positive.

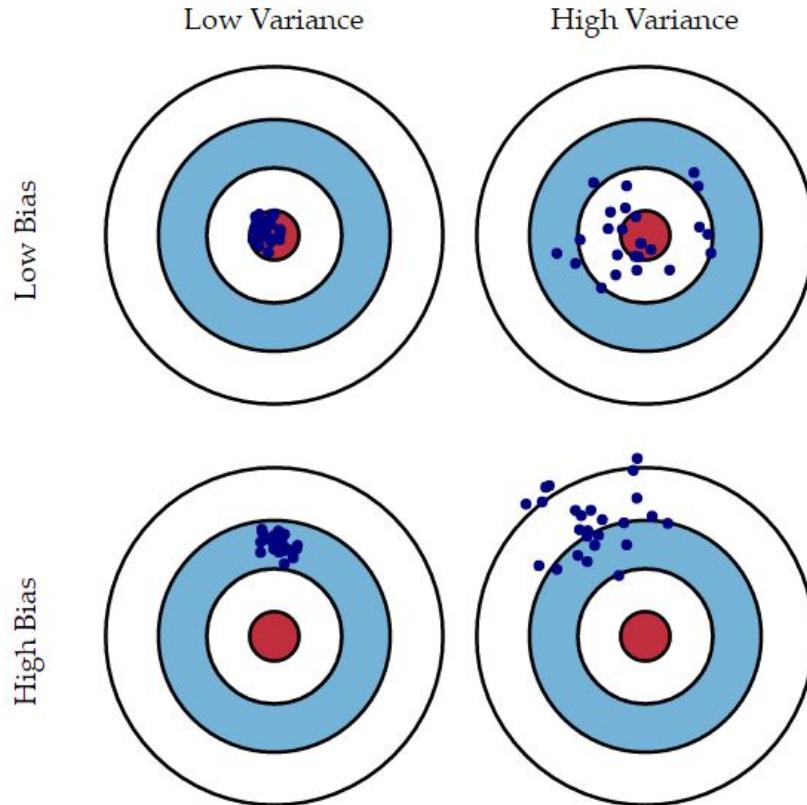
AUC → Metric to calculate the **overall performance** of a classification model based on area under the ROC curve.



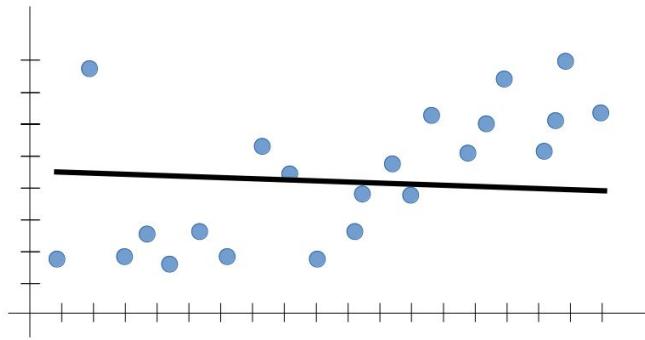
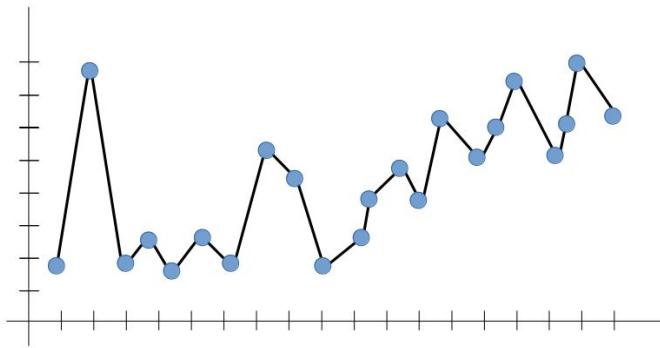
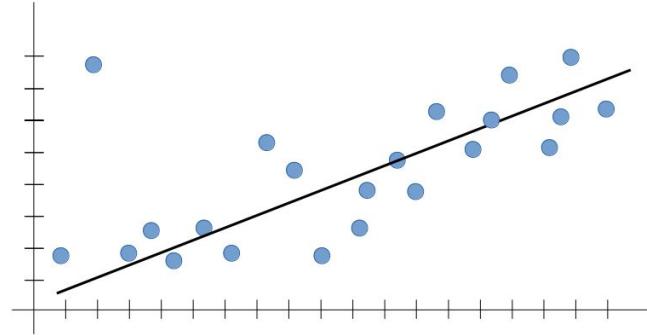
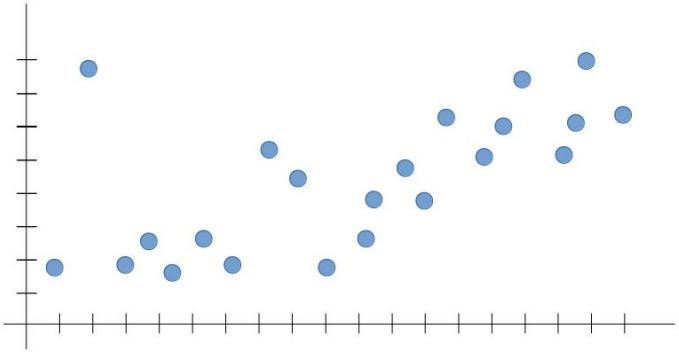
Cross-validation



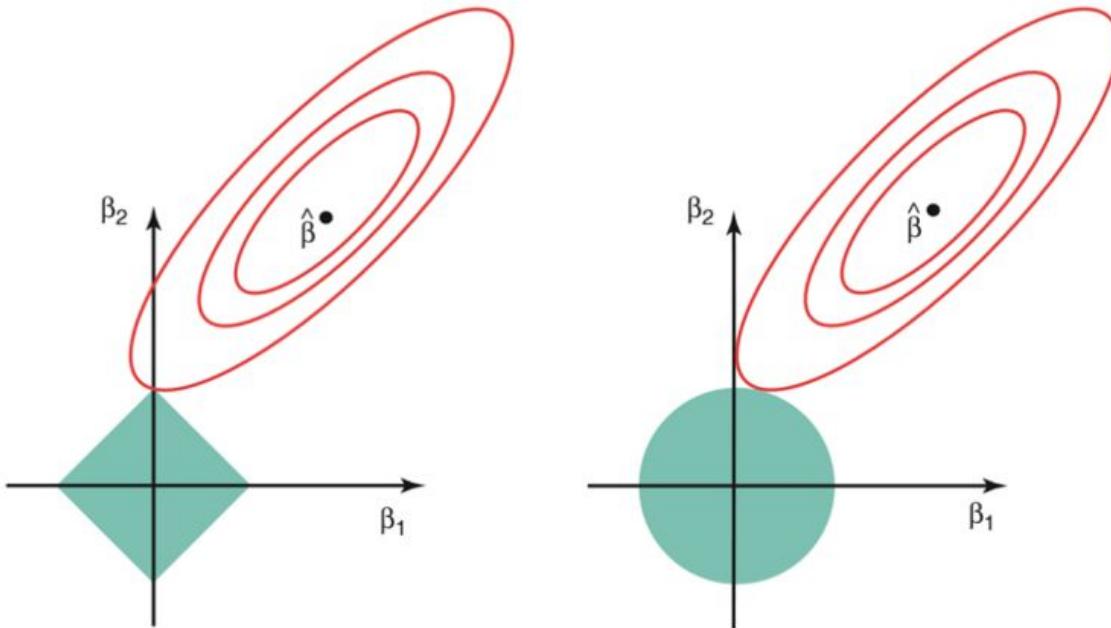
Bias vs. Variance



Underfitting vs. Overfitting



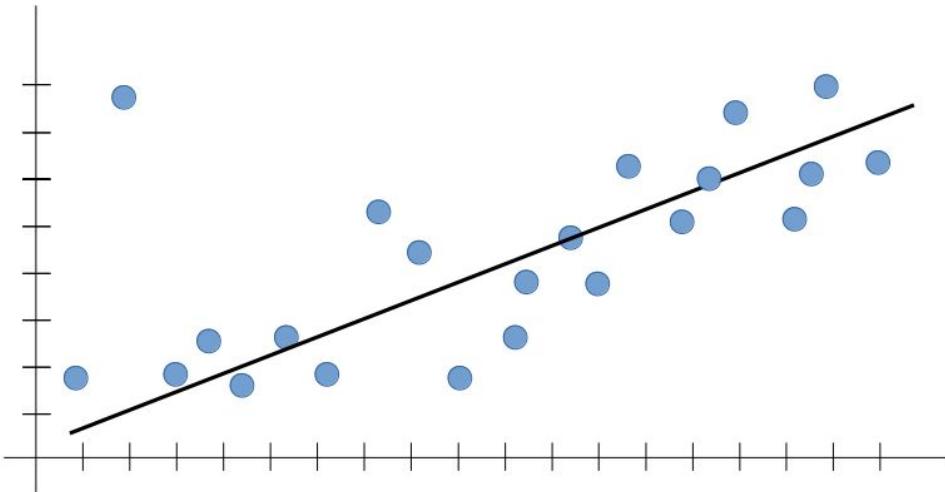
Regularization



A close-up, low-angle shot of the back and side of a dark, futuristic Iron Man suit. The suit has glowing red accents on the shoulder and knee areas. The character is positioned in front of a large, illuminated control panel with multiple screens displaying complex data and code. The overall atmosphere is dark and high-tech.

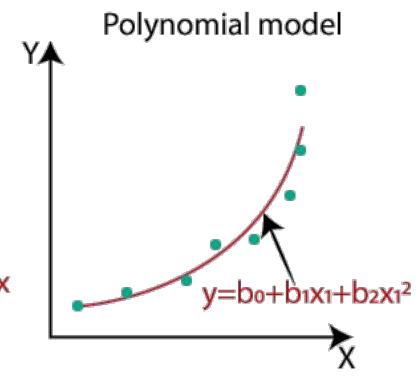
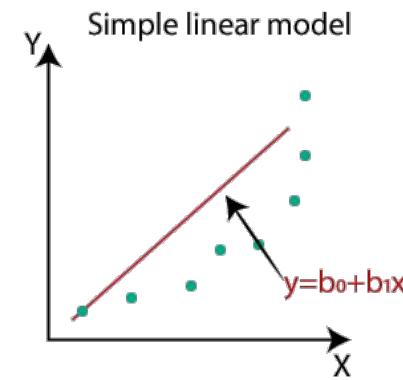
Let's dive into models! 😎

Linear Regression



Entrada: Dos dimensiones, datos
Salida: Un valor estimado (número)

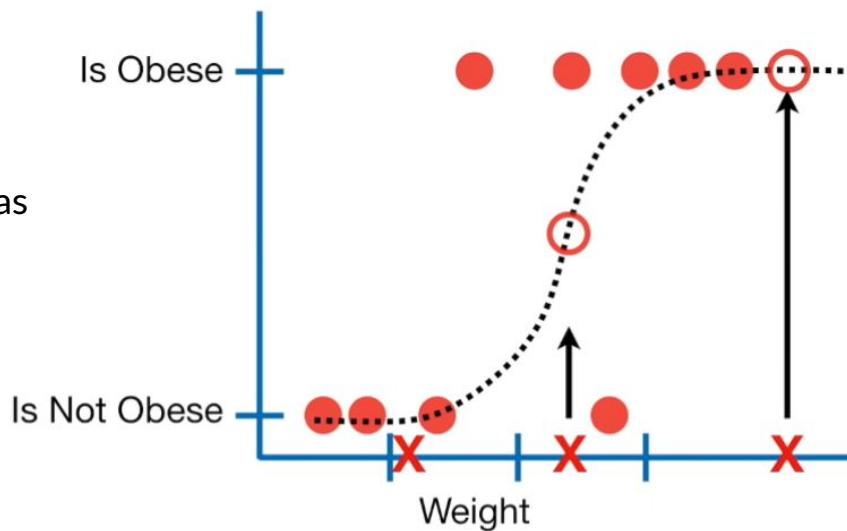
¿Correlación?



Logistic Regression

For example, if the probability a mouse is obese is $> 50\%$, then we'll classify it as obese, otherwise we'll classify it as "not obese".

Entrada: Una dimensión, dos categorías
Salida: Probabilidad (**clasificación**)



Naive Bayes

Algoritmo probabilístico. Teoría → Conocimiento parcial, proceso estocástico (no determinista)

Key: variables → tabla de frecuencia → probabilidades e incertidumbre

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

Probability of A occurring given evidence B has already occurred

Probability of B occurring given evidence A has already occurred

Probability of A occurring

Probability of B occurring

Naive Bayes

Dataset

| Weather | Play |
|----------|------|
| Sunny | No |
| Overcast | Yes |
| Rainy | Yes |
| Sunny | Yes |
| Sunny | Yes |
| Overcast | Yes |
| Rainy | No |
| Rainy | No |
| Sunny | Yes |
| Rainy | Yes |
| Sunny | No |
| Overcast | Yes |
| Overcast | Yes |
| Rainy | No |

Tabla de Frecuencias

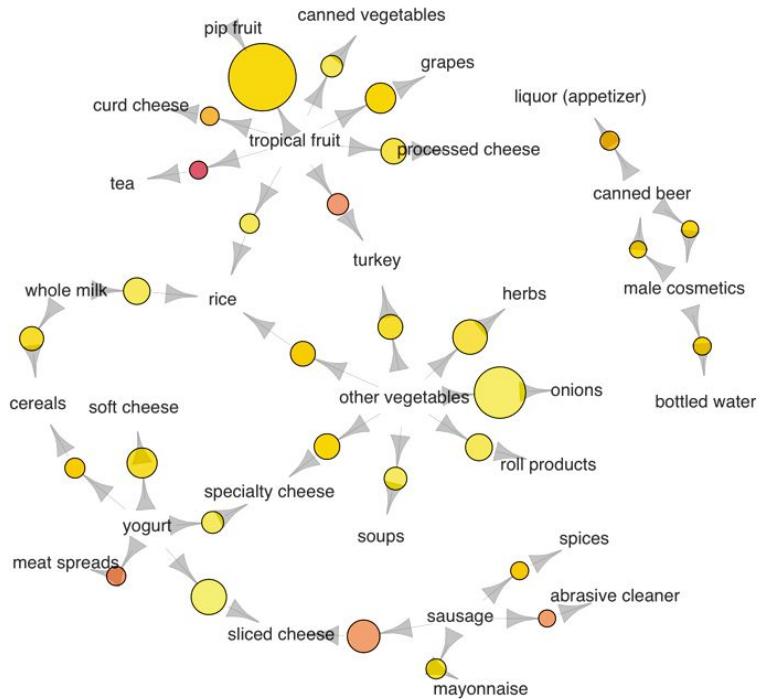
| Frequency Table | | |
|-----------------|----|-----|
| Weather | No | Yes |
| Overcast | | 4 |
| Rainy | 3 | 2 |
| Sunny | 2 | 3 |
| Grand Total | 5 | 9 |

Probabilidades

| Likelihood table | | | | |
|------------------|----|-------|-------|------|
| Weather | No | Yes | | |
| Overcast | | 4 | =4/14 | 0.29 |
| Rainy | 3 | 2 | =5/14 | 0.36 |
| Sunny | 2 | 3 | =5/14 | 0.36 |
| All | 5 | 9 | | |
| | | =5/14 | =9/14 | |
| | | 0.36 | 0.64 | |

→ Laplace trick

Association Rules



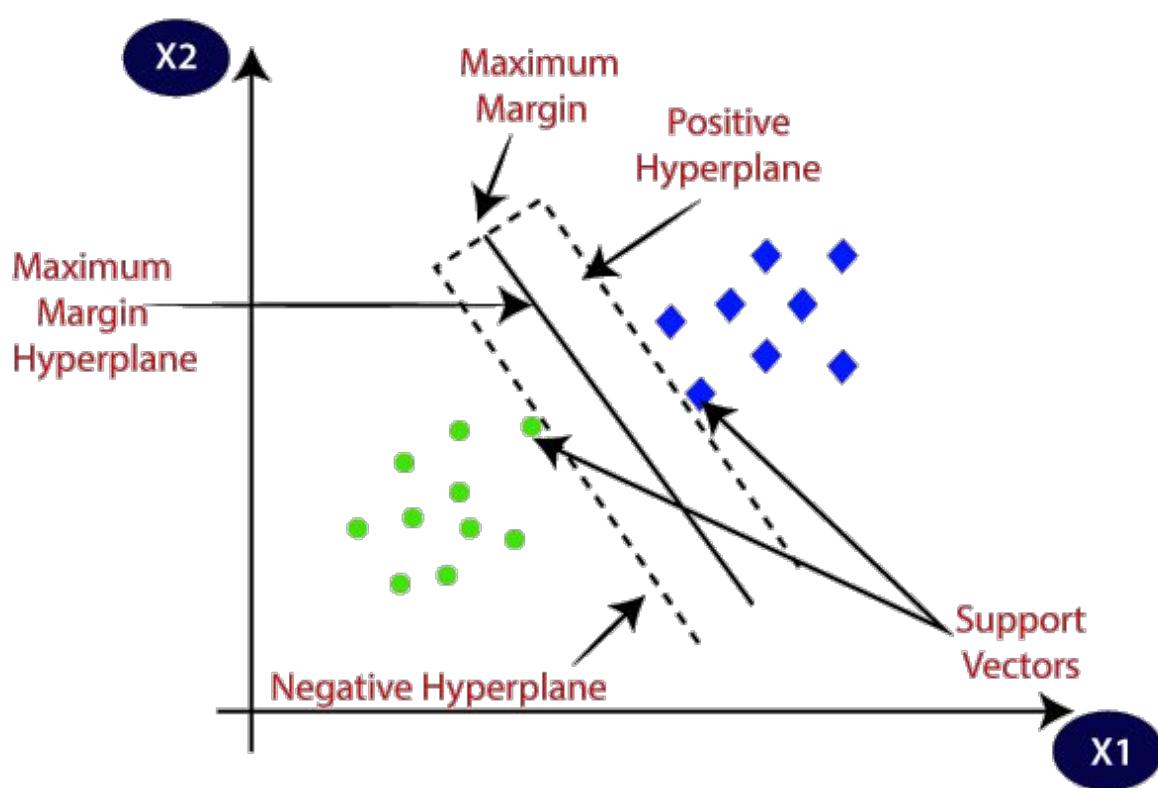
$$Rule: X \Rightarrow Y$$

$Support = \frac{frq(X, Y)}{N}$
 $Confidence = \frac{frq(X, Y)}{frq(X)}$
 $Lift = \frac{Support}{Supp(X) \times Supp(Y)}$

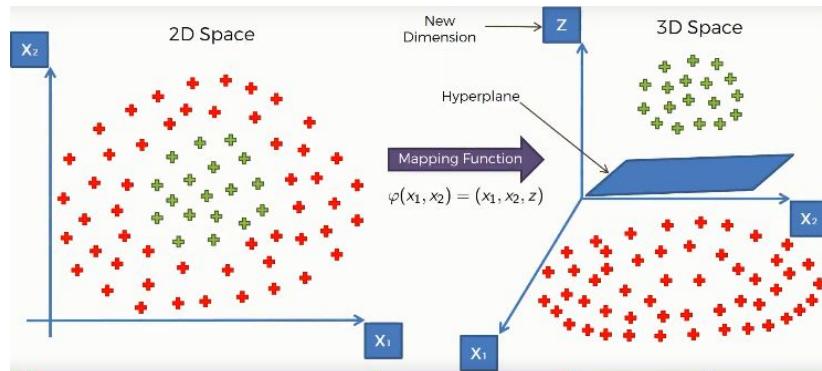
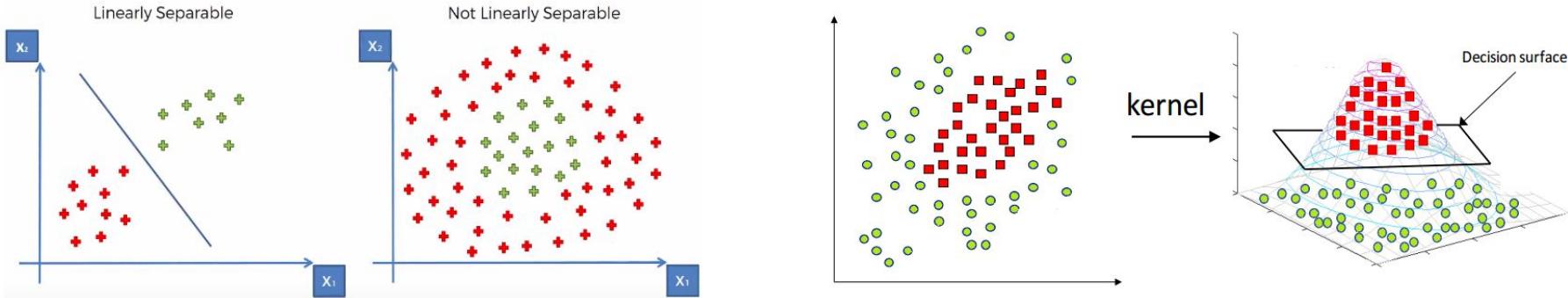
Association Rules

| | | | | |
|---------------|---|---|---|---|
| Transaction 1 | 🍎 | 🍺 | 🥣 | 🍗 |
| Transaction 2 | 🍎 | 🍺 | 🥣 | |
| Transaction 3 | 🍎 | 🍺 | | |
| Transaction 4 | 🍎 | 🍐 | | |
| Transaction 5 | 🍼 | 🍺 | 🥣 | 🍗 |
| Transaction 6 | 🍼 | 🍺 | 🥣 | |
| Transaction 7 | 🍼 | 🍺 | | |
| Transaction 8 | 🍼 | 🍐 | | |

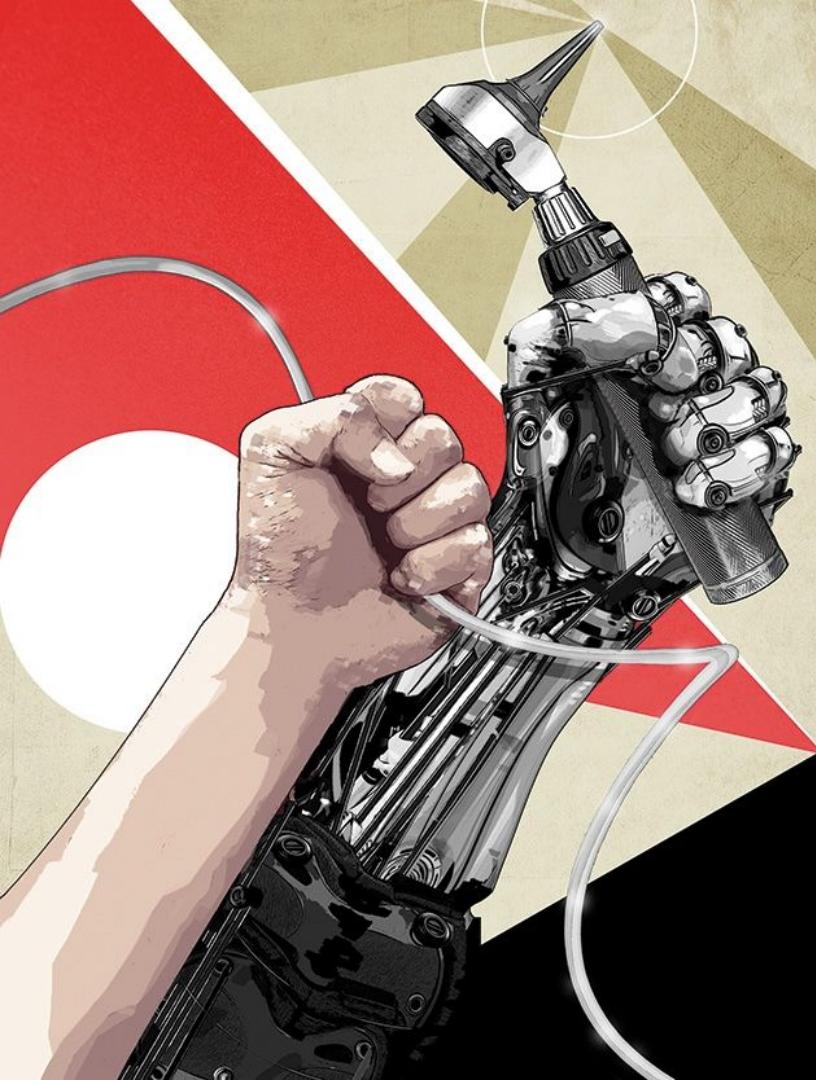
Support Vector Machine



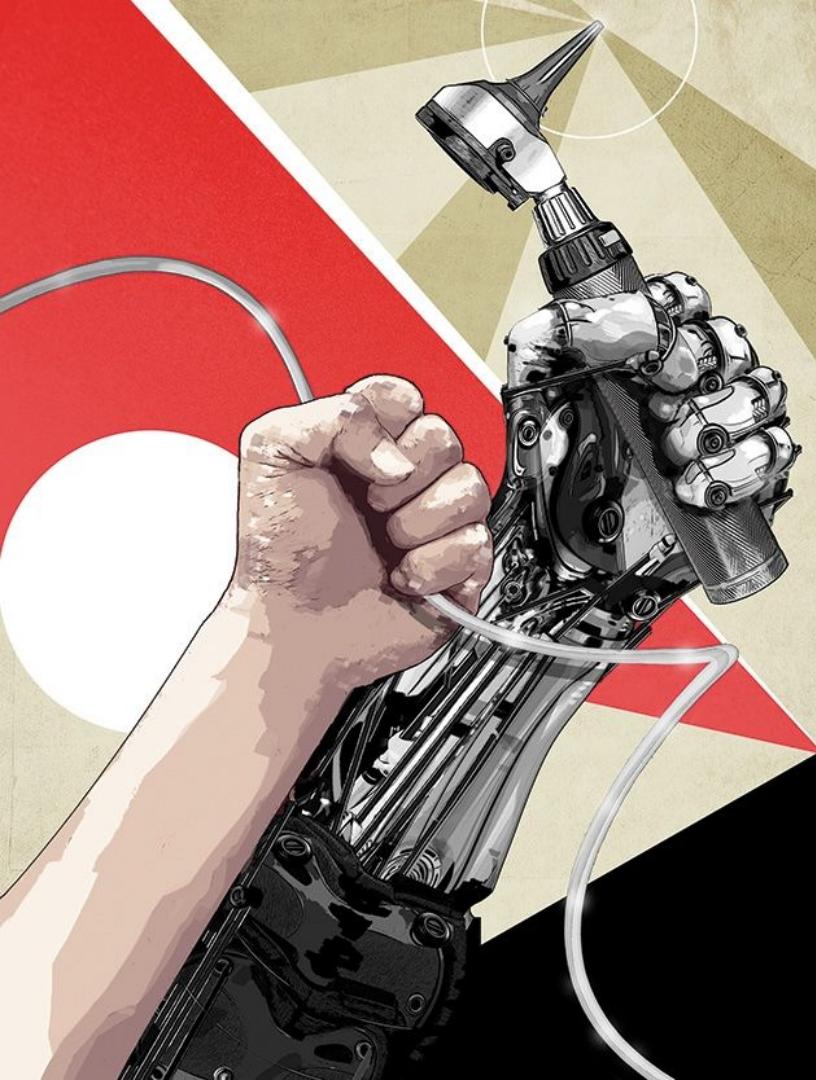
Support Vector Machine



Practice - Supervised Learning!



Challenge - Supervised Learning!



Bibliografía

/1./ /Hands-on Machine Learning with Scikit-Learn, Keras, and TensorFlow/

/2./ /Fast.AI - Introduction to Machine Learning for Coders/

/3./ /MLCourse.AI/

/4./ /DeltaAnalytics/

/5./ /The Hundred-page Machine Learning Book/

/6./ /Machine Learning for Humans (Vishal Maini)/

/7./ /Datacamp/

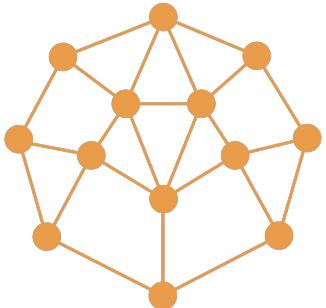
/8./ /DataQuest/



Partners

Agradecemos a nuestros partners por confiar en **nosotros** para facilitar la formación en IA de cara a la 4^a Revolución Industrial.



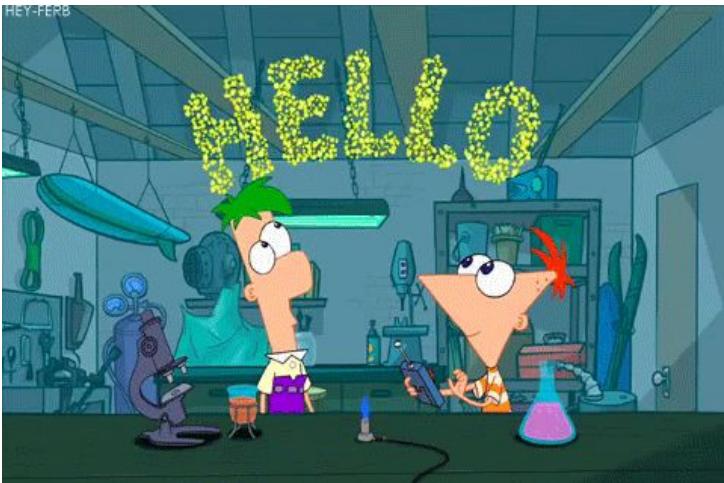


Saturdays.AI

This model fits me
95% of the time



WELCOME!



www.saturdays.ai
donostia.saturdays.ai
donostia@saturdays.ai